

Prevalence-Incidence Mixture Models

Description

This package fits Prevalence Incidence Mixture models to data for which the time to event is interval censored or for which the event is prevalent at the time zero but is only partially observed. Such data often arises in medical screening for asymptomatic disease or disease precursors, such as precancerous lesions. In such data, 1) onset of incident disease occurs between screening visits (interval-censoring), 2) the disease may have already occurred before the initial screen (prevalent at time zero) but may be initially missed and found some time after the initial screen. These models estimates absolute and relative risks. Semi-parametric (iterative convex minorant algorithm by Robertson, Wright and Dykstra, 1988), weakly-parametric (integrated B-splines), and fully parametric members of the Prevalence-Incidence Mixture model family are supported. A non-parametric estimator (Turnbull 1976) is provided and is useful for checking parametric assumptions of fully parametric Prevalence Incidence Mixture models. Only weakly-parametric and semi-parametric models (no variance calculation) currently support stratified random samples in the two viewpoints, a superpopulation and a finite population. The superpopulation is to view the finite population of interest as an independent sample of size N from an infinite superpopulation. A later version will add this functionality for the logistic-Weibull and logistic-exponential models. Semi-parametric, weakly-parametric models, logistic-Weibull, and logistic-exponential uses a logistic regression model as the prevalence model and a proportion hazard survival model as the incidence model. The semi-parametric model makes no assumptions regarding the baseline hazard function. However, it can be computationally expensive when the unique number of visit times are over hundreds. The weakly-parametric model approximates the baseline hazard function using integrated B-splines and is faster. When parametric assumptions can be made, the fully parametric models are fastest. The following parametric assumptions are supported: logistic-Weibull, logistic-exponential, logistic-lognormal, logistic-loglogistic, logistic-gengamma, and logistic-gamma. Variance estimates are available only for the weakly-parametric, logistic-Weibull, logistic-exponential, and non-parametric. For the non-parametric, this is achieved through boot-strapping by setting the `conf.int` parameter to "TRUE" and can be very computationally expensive. For identifiability of the mixture model, the data must contained observed prevalent disease and interval-censored incident disease.

Usage

```
PIixture(p.model, i.model, data, model = "semi-parametric",
  reg.initials = NULL, conf.int = FALSE,
  convergence.criteria = 0.001, iteration.limit = 250,
  time.interval = 0.01, design.out = TRUE, sample.design = NULL,
  N = NULL, n.knots = 5, order = 4, max.time, ...)
```

Arguments

- | | |
|----------------------|---|
| <code>p.model</code> | <p>The prevalence model to be fitted, specified using an expression of the form <i>c+l+r~model</i>. Elements in the expression are as follows:</p> <ul style="list-style-type: none"> • <code>c</code> - Numeric variable indicating whether the event was prevalent at time zero, taking values of 1=="Yes", 0=="No", -999="Latent"; • <code>l</code> - Numeric starting time of the interval in which event occurred, with -999 denoting known prevalent events; • <code>r</code> - Ending time of the interval in which event occurred, with -999 and Inf denoting known prevalent events and right-censoring, respectively; • <code>model</code> - Linear predictor consisting of a series of terms separated by <code>+</code> operators. |
| <code>i.model</code> | <p>The incidence model to be fitted, specified using an expression of the form <i>c+l+r~model</i> (see <code>p.model</code>). Defaults to <code>p.model</code>.</p> |
| <code>data</code> | <p>Data used to fit the model containing columns for each term in <code>p.model</code> and <code>i.model</code> expressions. For stratified random sampling designs, columns denoted <code>samp.weight</code> and <code>strata</code> are expected indicating the sampling weights and sampling strata. For <code>sample.design="superpopulation"</code> option, an additional column denoted <code>strata.frac</code> is expected indicating the fraction of the population that consists of each strata. For example, if in the target population there are three strata that occurs with proportions 0.2, 0.4, and 0.6, then <code>strata.frac</code> will take values of 0.2, 0.4 or 0.6.</p> |
| <code>model</code> | <p>Character string indicating the specific member of the Prevalence-Incidence Mixture Model family to be fitted. Options are:</p> <ul style="list-style-type: none"> • "semi-parametric" Fits logistic regression and proportional hazards model as the prevalence and incidence models, respectively. The baseline hazard function is non-parametrically estimated using the iterative convex minorant algorithm. • "weakly-parametric" Fits logistic regression and proportional hazards model as the prevalence and incidence models, respectively. The baseline hazard function is approximated using integrated B-splines. • "logistic-Weibull" Fits logistic regression and proportional hazards model as the prevalence and incidence models, respectively. The baseline hazard function is approximated using a Weibull distribution. • "logistic-exponential" Fits logistic regression and proportional hazards model as the prevalence and incidence models, respectively. The baseline hazard function is approximated using an exponential distribution. • "logistic-lognormal" Fits logistic regression and lognormal survival as the prevalence and incidence models, respectively. • "logistic-loglogistic" Fits logistic regression and loglogistic survival as the prevalence and incidence models, respectively. • "logistic-gengamma" Fits logistic regression and generalized-gamma survival as the prevalence and incidence models, respectively. • "logistic-gamma" Fits logistic regression and gamma survival as the prevalence and incidence models, respectively. • "non-parametric" Provides the non-parametric cumulative risk estimator. This is akin to the non-parametric estimates provided using the Turnbull methods. Covariates are not supported. Confidence intervals are obtained through bootstrapping but is computationally expensive. |

Variance estimates are not available for "semi-parametric", "logistic-lognormal", "logistic-loglogistic", "logistic-gengamma", or "logistic-gamma", but they can be obtained using bootstrap methods. Defaults to "semi-parametric".

<code>reg.initials</code>	Initial parameter estimates. Defaults to NULL.
<code>conf.int</code>	For non-parametric model option, FALSE="Do not obtain bootstrap confidence intervals", TRUE="Obtain bootstrap confidence intervals. Defaults to FALSE.
<code>convergence.criteria</code>	Convergence of models occurs when reduction in the objective is within this convergence tolerance. Defaults to 0.001.
<code>iteration.limit</code>	Maximum number of iterations allowed to achieve convergence. Defaults to 250.
<code>time.interval</code>	Define time intervals to output baseline hazards for. Defaults to .01.
<code>design.out</code>	Option to include the design matrix of data used for model fitting in the output. Defaults to TRUE.
<code>sample.design</code>	Sampling design of the NULL="simple random sampling", 1="finite population", 2="superpopulation". Defaults to NULL. For "superpopulation", N is required for variance calculation (only provided when model="weakly-parametric" is used)
<code>N</code>	Population size, required for superpopulation. Defaults to NULL.
<code>n.knots</code>	Number of knots for splines for "weakly-parametric" model. Defaults to 5.
<code>order</code>	Degree of splines for "weakly-parametric" model. Defaults to 4 (cubic splines).
<code>max.time</code>	Define maximum time to output baseline hazards for. Defaults to the largest finite start/end time given in the data.

Value

The output is a list of class PIMix which contains the following elements.

- `data.summary` A data frame containing the following: Num. of observation - total number of observations in data set; Included subjects - number of observations used in fitting model; Known prevalent cases - the number of events known to be prevalent at time zero; Interval censoring - the number of event times occurring in the interval ($L > 0, R < \text{Inf}$]; Left censoring - the number of event times known to occur by $R < \text{Inf}$, but can also have been prevalent at time zero; Right censoring - the number of observations right-censored with event time occurring in the interval ($L > 0, \text{Inf}$); Missing prevalent status - the number of observations where it is unknown whether the event was prevalent at time zero; Non-informative intervals - the number of observations with intervals (0,Inf) or [0,Inf) (denoting missing prevalent status).
- `regression.coef` A data frame summarizing parameter values, standard errors, and 95 percent confidence intervals.
- `OR` A data frame summary odds ratios, , standard errors, and 95 percent confidence intervals.
- `HR` A data frame summary hazard ratios, , standard errors, and 95 percent confidence intervals.
- `knots` If model="weakly-parametric" is specified, this is a numeric vector of starting time points for each exponential spline.
- `exp.spline.coeff` If model="weakly-parametric" is specified, this is a numeric vector of coefficients for each exponential spline.
- `cum.hazard` If model="semi-parametric" or model="weakly-parametric" is specified, this is a data frame containing the baseline cumulative hazard.
- `covariance` A matrix containing the covariance matrix for the parameters (not produced for model="semi-parametric").
- `hessian` A matrix containing the hessian matrix for the parameters (not produced for model="semi-parametric").
- `model` Character string indicating the specific member of the Prevalence-Incidence Mixture Model family fitted.
- `p.model` The prevalence model.
- `prev.design` The design matrix for the prevalence model.
- `i.model` The incidence model.
- `incid.design` The design matrix for the incidence model.
- `loglikelihood` For random samples, this is the log-likelihood of the fitted model. For stratified random samples, the weighted-likelihood approach is used and a log-pseudolikelihood (weighted log-likelihood) is reported.
- `convergence` Convergence statistics.

Warning

The model="semi-parametric" option is very computationally expensive when the unique visit times are over hundreds.

Author(s)

Li C. Cheung, li.cheung@nih.gov, Noorie Hyun nhyun@mcw.edu Xiaojin Xiong, Qing Pan, Hormuzd A. Katki

References

- Cheung LC, Qing P, Hyun N, Schiffman M, Fetterman B, Castle P, Lorey T, Katki H. Mixture models for undiagnosed prevalent disease and interval-censored incident disease: Applications to a cohort assembled from electronic health records. *Statistics in Medicine* 2017; 36(22):3583-95.
- Hyun N, Cheung LC, Pan Q, Katki H. Flexible risk prediction models for left or interval-censored data from electronic health records. *Annals of Applied Statistics* 11(2), 1063-1084.
- Turnbull BW (1976). The empirical distribution with arbitrary grouped censored and truncated data. *Journal of the Royal Statistical Society - Series B (Statistical Methodology)* 38, 290-295.

- Robertson T, Wright FT, and Dykstra RL (1988). Order Restricted Statistical Inference. Wiley.

Examples

```
#PIMixture includes "PIdata" RData file, and PIdata includes the two datasets, PIdata1 and PIdata2
data(PIdata)
model<-"C_CIN3PLUS+L_CIN3PLUS+R_CIN3PLUS~RES_HPV16"
fit1<-PIMixture(p.model=model,data=PIdata1, model="logistic-Weibull")
fit2<-PIMixture(p.model=model,data=PIdata1, model="weakly-parametric",n.knots=5,order=4)
fit3<-PIMixture(p.model=model,data=PIdata1, model="semi-parametric")

model2<-"C_CIN3PLUS+L_CIN3PLUS+R_CIN3PLUS~1"
fit4<-PIMixture(p.model=model2,data=PIdata1, model="non-parametric", conf.int=TRUE)

#For stratified random samples
model3<-"C+L+R~X1+X2"
#sample.design=1 indicates the target population is a finite population, and the variance is design-based.
fit5<-PIMixture(p.model=model3,data=PIdata2, model="weakly-parametric",n.knots=7,order=4,sample.design=1)

#sample.design=2 indicates the target population is a superpopulation, and the variance consists of
#design-based and model-based variances. Generally, the Variance in the superpopulation frame is slightly larger than
fit6<-PIMixture(p.model=model3,data=PIdata2, model="weakly-parametric",n.knots=7,order=4,sample.design=2,N=10000)

fit7<-PIMixture(p.model=model3,data=PIdata2, model="semi-parametric",sample.design=1)
```